# Volume II, Appendix C
# Table of Contents

**C**  ₁ # Appendix C: Qualification Test
₂ # Design Criteria

₃

₄ ## C.1      Scope

₅ This appendix describes the guiding principles used to design the voting system
₆ qualification testing process conducted by ITAs.

₇ Qualification tests are designed to demonstrate that the system meets or exceeds the
₈ requirements of the Standards. The tests are also used to demonstrate compliance with
₉ other levels of performance claimed by the manufacturer.

₁₀ Qualification tests must satisfy two separate and possibly conflicting sets of
₁₁ considerations. The first is the need to produce enough test data to provide confidence
₁₂ in the validity of the test and its apparent outcome. The second is the need to achieve a
₁₃ meaningful test at a reasonable cost, and cost varies with the difficulty of simulating
₁₄ expected real-world operating conditions and with test duration. It is the test
₁₅ designer's job to achieve an acceptable balance of these constraints.

₁₆ The rationale and statistical methods of the test designs contained in the Standards are
₁₇ discussed below. Technical descriptions of their design can be found in any of several
₁₈ books on testing and statistical analysis.

₁₉ ## C.2      Approach to Test Design

₂₀ The qualification tests specified in the Standards are primarily concerned with
₂₁ assessing the magnitude of random errors. They are also, however, capable of
₂₂ detecting bias errors that would result in the rejection of the system.

₂₃ Test data typically produce two results. The first is an estimate of the true value of
₂₄ some system attribute such as speed, error rate, etc. The second is the degree of
₂₅ certainty that the estimate is a correct one. The estimate of an attribute's value may or
₂₆ may not be greatly affected by the duration of the test. Test duration, however, is very

1    important to the degree of certainty; as the length of the test increases, the level of
2    uncertainty decreases. An efficient test design will produce enough data over a
3    sufficient period of time to enable an estimate at the desired level of confidence.

4    There are several ways to design tests. One approach involves the preselection of
5    some test parameter, such as the number of failures or other detectable factor. The
6    essential element of this type of design is that the number of observations is
7    independent of their results. The test may be designed to terminate after 1,000 hours
8    or 10 days, or when 5 failures have been observed. The number of failures is
9    important because the confidence interval (uncertainty band) decreases rapidly as the
10   number of failures increases. However, if the system is highly reliable or very
11   accurate, the length of time required to produce a predetermined number of failures or
12   errors using this method may be unachievably long.

13   Another approach is to determine that the actual value of some attribute need not be
14   learned by testing, provided that the value can be shown to be better than some level.
15   The test would not be designed to produce an estimate of the true value of the attribute
16   but instead to show, for example, that reliability is at least 123 hours or the error rate
17   is no greater than one in ten million characters.

18   The latter design approach, which was chosen for the Standards, uses what is called
19   Sequential Analysis. Instead of the test duration being fixed, it varies depending on
20   the outcome of a series of observations. The test is terminated as soon as a statistically
21   valid decision can be reached that the factor being tested is at least as good as or no
22   worse than the predetermined target value. A sequential analysis test design called the
23   "Wald Probability Ratio Test" is used for reliability and accuracy testing.

## 24   C.3      Probability Ratio Sequential Test (PRST)

25   The design of a Probability Ratio Sequential Test (PRST) requires that four
26   parameters be specified:

27          H0, the null hypothesis
28          H1, the alternate hypothesis

29          a, the Producer's risk
30          b, the Consumer's risk

31   The Standards anticipate using the PRST for testing both time-based and event-based
32   failures.

33   This test design provides decision criteria for accepting or rejecting one of two test
34   hypotheses:  the null hypothesis, which is the Nominal Specification Value (NSV), or
35   the alternate hypothesis, which is the MAV. The MAV could be either the Minimum
36   Acceptable Value or the Maximum Acceptable Value depending upon what is being

1    tested. (Performance may be specified by means of a single value or by two values.
2    When a single value is specified, it shall be interpreted as an upper or lower single-
3    sided 90 percent confidence limit. If two values, these shall be interpreted as a two-
4    sided 90 percent confidence interval, consisting of the NSV and MAV.)

5    In the case of Mean Time Between Failure (MTBF), for example, the null hypothesis
6    is that the true MTBF is at least as great as the desired value (NSV), while the
7    alternate hypothesis is that the true value of the MTBF is less than some lower value
8    (Minimum Acceptable Value). In the case of error rate, the null hypothesis is that the
9    true error rate is less than some very small desired value (NSV), while the alternate
10   hypothesis is that the true error rate is greater than some larger value that is the upper
11   limit for acceptable error (Maximum Acceptable Value).

12   ## C.4    Time-based Failure Testing Criteria

13   An equivalence between a number of events and a time period can be established
14   when the operating scenarios of a system can be determined with precision.  Some of
15   the performance test criteria of Volume II, Section 4, *Hardware Testing,* use this
16   equivalence.

17   System acceptance or rejection can be determined by observing the number of
18   relevant failures that occur during equipment operation. The probability ratio for this
19   test is derived from the Exponential probability distribution. This distribution implies
20   a constant hazard rate for equipment failure that is not dependent on the time of
21   testing or the previous failures.  In that case, two or more systems may be tested
22   simultaneously to accumulate the required number of test hours, and the validity of
23   the data is not affected by the number of operating hours on a particular unit of
24   equipment. However, for environmental operating hardware tests, no unit shall be
25   subjected to less than two complete 24 hour test cycles in a test chamber as required
26   by Volume II, Subsection 4.7.1 of the Standards.

27   In this case, the null hypothesis is that the Mean Time Between Failure (MTBF), as
28   defined in Volume I, Subsection 3.4.3 of the Standards, is at least as great as some
29   value, here the Nominal Specification Value. The alternate hypothesis is that the
30   MTBF is no better than some value, here the Minimum Acceptable Value.

31   For example, a typical system operations scenario for environmental operating
32   hardware tests will consist of approximately 45 hours of equipment operation. Broken
33   down, this time allotment involves 30 hours of equipment set-up and readiness testing
34   and 15 hours of elections operations. If the Minimum Acceptable Value is defined as
35   45 hours, and a test discrimination ratio of 3 is used (in order to produce an acceptably
36   short expected time of decision), then the Nominal Specification Value equals 135
37   hours.

1  With a value of decision risk equal to 10 percent, there is no more than a 10 percent
2  chance that a system would be rejected when, in fact, with a true MTBF of at least 135
3  hours, the system would be acceptable. It also means that there is no more than a 10
4  percent chance that a system would be accepted with a true MTBF lower than 45
5  hours when it should have been rejected.

6  Therefore,

7  H0:  MTBF = 135 hours
8  H1:  MTBF = 45 hours

9                              a =      0.10
10                             b =      0.10.

11  Under this PRST design, the test is terminated and an ACCEPT decision is reached
12  when the cumulative number of equipment hours in the second column of the
13  following table has been reached, and the number of failures is equal to or less than
14  the number shown in the first column. The test is terminated and a REJECT decision
15  is reached when the number of failures occurs in less than the number of hours
16  specified in the third column.  Here, the minimum time to accept (on zero failures) is
17  169 hours.  In the event that no decision has been reached by the times shown in the
18  last table entries, the test is terminated, and the decision is declared as indicated. Any
19  time that 7 or more failures occur, the test is terminated and the equipment rejected.
20  If after 466 hours of operation the cumulative failure score is less than 7.0, then the
21  equipment is accepted.

22

| 23 Number of | Accept if Time | Reject if Time |
|---|---|---|
| 24 Failures | Greater Than | Less Than |
| 25  0 | 169 | Continue test |
| 26  1 | 243 | Continue test |
| 27  2 | 317 | 26 |
| 28  3 | 392 | 100 |
| 29  4 | 466 | 175 |
| 30  5 | 466 | 249 |
| 31  6 | 466 | 323 |
| 32  7 | N/A | (1) |

33                        (1) Terminate and REJECT

34

1 This test is based on the table of test times of the truncated PRST design V-D in the
2 Military Handbook MIL-HDBK-781A that is designated for discrimination ratio 3 and
3 a nominal value of 0.10 for both a and b.  The Handbook states that the true producer
4 risk is 0.111 and the true consumer risk is 0.109.   Using the theoretical formulas for
5 either the untruncated or truncated tests will lead to different numbers.

6 The test design will change if given a different set of parameters.   Some jurisdictions
7 may find the Minimum Acceptable Value of 45 hours unacceptable for their needs.  In
8 addition, it may be appropriate to use a different discrimination ratio, or different
9 Consumer's and Producer's risk.  Also, before using tests based on the MTBF, it
10 should be determined whether time-based testing is appropriate rather than event-
11 based or another form of testing.  If MTBF-based procedures are chosen, then the
12 appropriateness of the assumption of a constant hazard rate with exponential failures
13 should in turn be assessed.

14

15 # C.5    Accuracy Testing Criteria

16 Some voting system performance attributes are tested by inducing an event or series
17 of events, and the relative or absolute time intervals between repetitions of the event
18 has no significance. Although an equivalence between a number of events and a time
19 period can be established when the operating scenarios of a system can be determined
20 with precision, another type of test is required when such equivalence cannot be
21 established. It uses event-based failure frequencies to arrive at ACCEPT/REJECT
22 criteria. This test may be performed simultaneously with time-based tests.

23 For example, the failure of a device is usually dependent on the processing volume
24 that it is required to perform. The elapsed time over which a certain number of
25 actuation cycles occur is, under most circumstances, not important. Another example
26 of such an attribute is the frequency of errors in reading, recording, and processing
27 vote data.

28 The error frequency, called "ballot position error rate," applies to such functions as
29 process of detecting the presence or absence of a voting punch or mark, or to the
30 closure of a switch corresponding to the selection of a candidate.

31 Qualification and acceptance test procedures that accommodate event-based failures
32 are, therefore, based on a discrete, rather than a continuous probability distribution. A
33 Probability Ratio Sequential Test using the binomial distribution is recommended. In
34 the case of ballot position error rate, the calculation for a specific device (and the
35 processing function that relies on that device) is based on:

36         HO: Desired error rate = 1 in 10,000,000

1    H1: Maximum acceptable error rate = 1 in 500,000

2    a = 0.05

3    b= 0.05

4    and the minimum error-free sample size to accept for qualification tests is 1,549,703
5    votes.

6    The nature of the problem may be illustrated by the following example, using the
7    criteria contained in the Standards for system error rate. A target for the desired
8    accuracy is established at a very low error rate. A threshold for the worst error rate
9    that can be accepted is then fixed at a somewhat higher error rate. Next, the decision
10   risk is chosen, that is the risk that the test results may not be a true indicator of either
11   the system's acceptability or unacceptability. The process is as follows:

12   ♦   The desired accuracy of the voting system, whatever its true error rate (which
13       may be far better), is established as no more than one error in every ten
14       million characters (including the null character).

15   ♦   If it can be shown that the system's true error rate does not exceed one in
16       every five hundred thousand votes counted, it will be considered acceptable.
17       (This is more than accurate enough to declare the winner correctly in almost
18       every election.)

19   ♦   A decision risk of 5 percent is chosen, to be 95 percent sure that the test data
20       will not indicate that the system is bad when it is good or good when it is bad.

21   This results in the following decision criteria:

22   ♦   If the system makes one error before counting 26,997 consecutive ballot
23       positions correctly, it will be rejected. The vendor is then required to improve
24       the system;

25   ♦   If the system reads at least 1,549,703 consecutive ballot positions correctly, it
26       will be accepted; and

27   If the system correctly reads more than 26,997 ballot positions but less than
28       1,549,703 when the first error occurs, the testing will have to be continued
29       until another 1,576,701 consecutive ballot positions are counted without error
30       (a total of 3,126,404 with one error).